

Linear Regression

Marc Deisenroth

Centre for Artificial Intelligence

Department of Computer Science

University College London

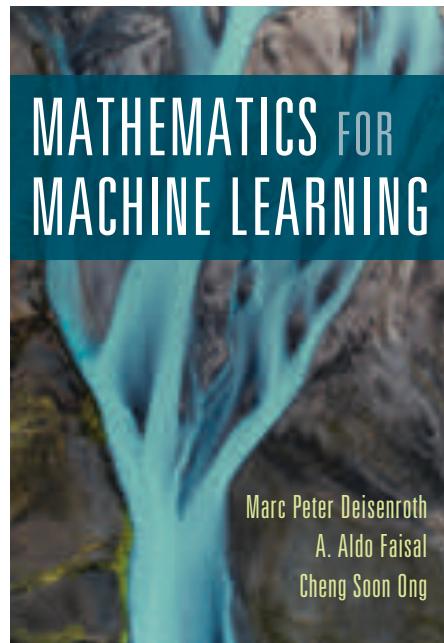
AIMS Rwanda and AIMS Ghana

March/April 2020

 @mpd37

m.deisenroth@ucl.ac.uk

<https://deisenroth.cc>

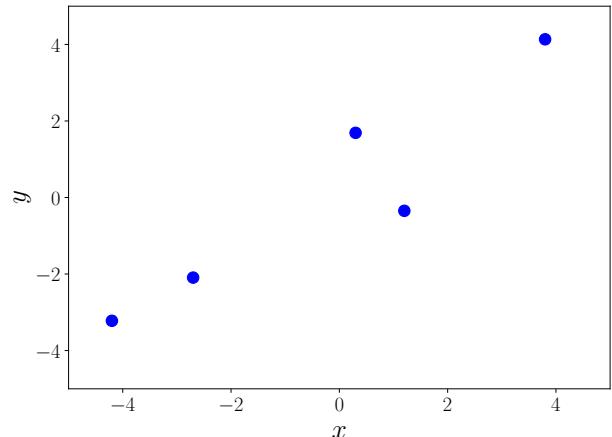


<https://mml-book.com>

Chapter 9

Regression (curve fitting)

Given inputs $x \in \mathbb{R}^D$ and corresponding observations $y \in \mathbb{R}$ find a function f that models the relationship between x and y .

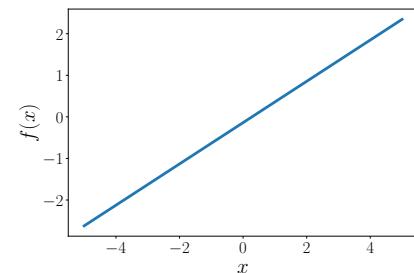


- Typically parametrize the function f with parameters θ
- Linear regression: Consider functions f that are **linear in the parameters**

Linear Regression Functions

■ Straight lines

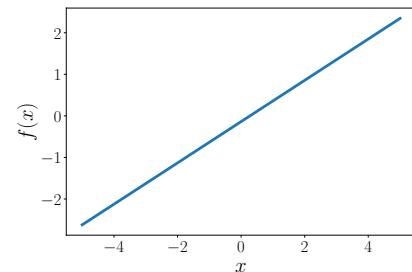
$$y = f(x, \theta) = \theta_0 + \theta_1 x = \begin{bmatrix} \theta_0 & \theta_1 \end{bmatrix} \begin{bmatrix} 1 \\ x \end{bmatrix}$$



Linear Regression Functions

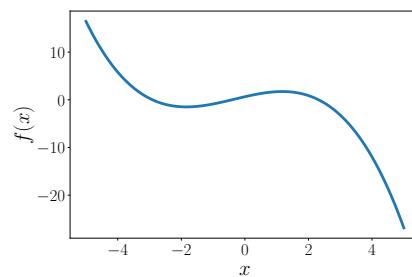
■ Straight lines

$$y = f(x, \theta) = \theta_0 + \theta_1 x = \begin{bmatrix} \theta_0 & \theta_1 \end{bmatrix} \begin{bmatrix} 1 \\ x \end{bmatrix}$$



■ Polynomials

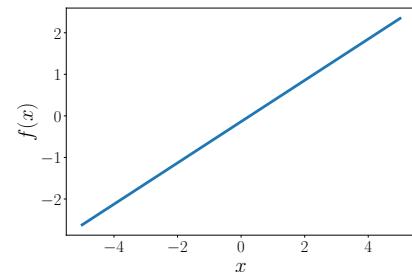
$$y = f(x, \theta) = \sum_{m=0}^M \theta_m x^m = \begin{bmatrix} \theta_0 & \dots & \theta_M \end{bmatrix} \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix}$$



Linear Regression Functions

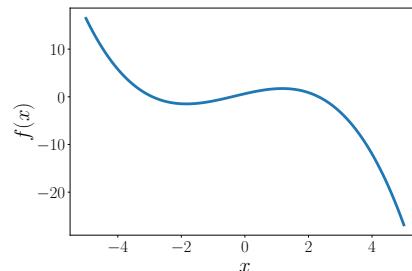
■ Straight lines

$$y = f(x, \theta) = \theta_0 + \theta_1 x = \begin{bmatrix} \theta_0 & \theta_1 \end{bmatrix} \begin{bmatrix} 1 \\ x \end{bmatrix}$$



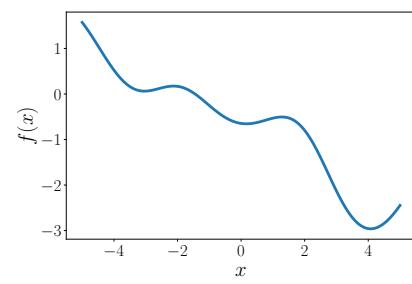
■ Polynomials

$$y = f(x, \theta) = \sum_{m=0}^M \theta_m x^m = \begin{bmatrix} \theta_0 & \dots & \theta_M \end{bmatrix} \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix}$$



■ Radial basis function networks

$$y = f(x, \theta) = \sum_{m=1}^M \theta_m \exp\left(-\frac{1}{2}(x - \mu_m)^2\right)$$

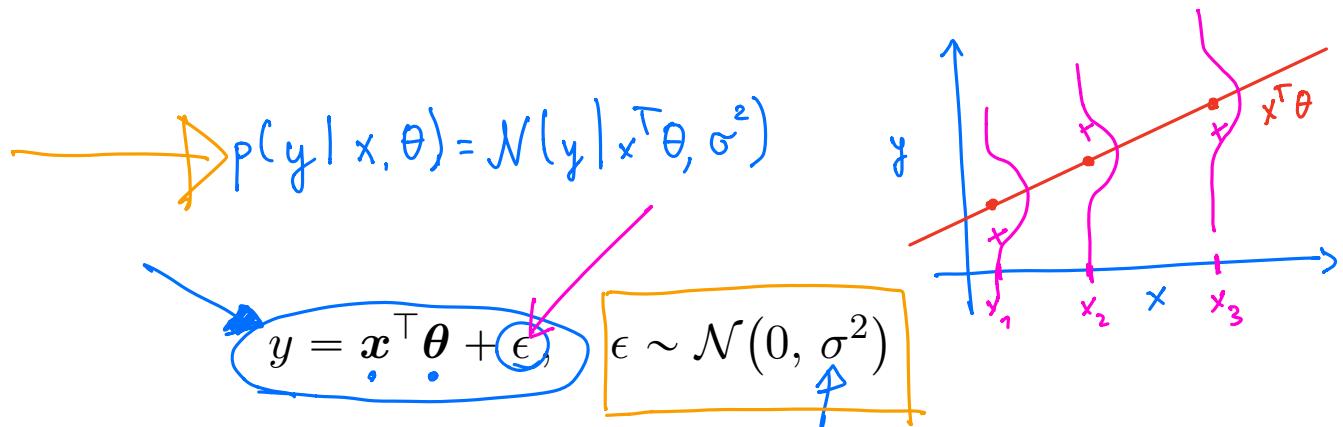


Linear Regression Model and Setting

$$y = \mathbf{x}^\top \boldsymbol{\theta} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

- Given a training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ we seek optimal parameters $\boldsymbol{\theta}^*$

Linear Regression Model and Setting



- Given a training set $(x_1, y_1), \dots, (x_N, y_N)$ we seek optimal parameters θ^*

► Maximum Likelihood Estimation

► Maximum a Posteriori Estimation

$$p(\varepsilon) = \mathcal{N}(0, \sigma^2)$$

$$\alpha := x^T \theta$$

$$y := \varepsilon + \alpha$$

$$\phi(\lambda x + \psi y) = \lambda \phi(x) + \psi \phi(y)$$

$$\text{var}(\lambda x) = \lambda^2 \text{var}(x)$$

$$\varepsilon + \alpha \Rightarrow p(\varepsilon + \alpha) = \mathcal{N}(m, s^2)$$

$$E[\varepsilon + \alpha] = \underbrace{E[\varepsilon]}_{=0} + \alpha = \alpha$$

$$\text{var}[\varepsilon + \alpha] = \text{var}[\varepsilon] = \sigma^2$$

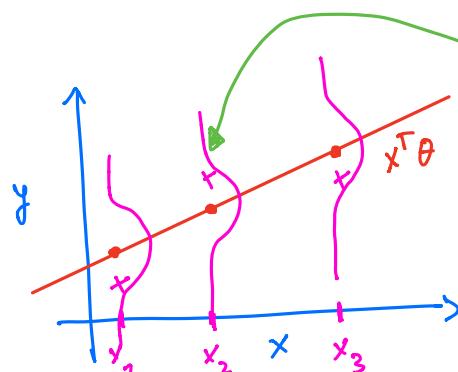
Maximum Likelihood

- Define $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D}$ and $\mathbf{y} = [y_1, \dots, y_N]^\top \in \mathbb{R}^N$
- Find parameters θ^* that maximize the likelihood

Maximum Likelihood

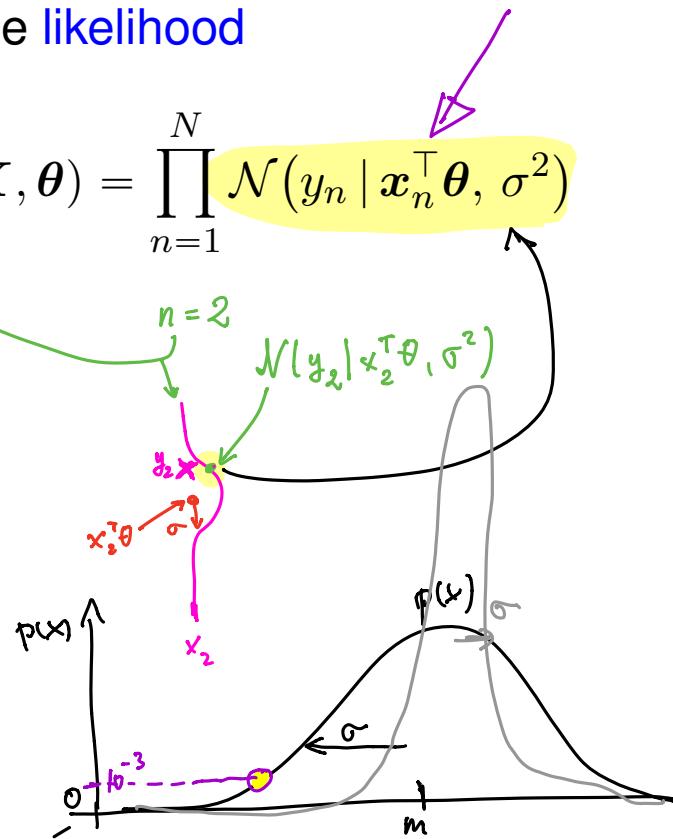
- Define $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D}$ and $\mathbf{y} = [y_1, \dots, y_N]^\top \in \mathbb{R}^N$
- Find parameters θ^* that maximize the likelihood

$$p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N, \boldsymbol{\theta}) = p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2)$$



$$\int p(x) dx = 1$$

$$\frac{N}{\pi} 10^{-3} = 10^{-3N}$$



- Define $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D}$ and $\mathbf{y} = [y_1, \dots, y_N]^\top \in \mathbb{R}^N$
- Find parameters θ^* that maximize the likelihood

$$p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N, \boldsymbol{\theta}) = p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2)$$

- Log-transformation ➤ **Maximize the log likelihood**

Maximum Likelihood

- Define $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D}$ and $\mathbf{y} = [y_1, \dots, y_N]^\top \in \mathbb{R}^N$
- Find parameters θ^* that maximize the likelihood

$$p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N, \boldsymbol{\theta}) = p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2)$$

- Log-transformation ➤ Maximize the log likelihood

$$\log p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) = \sum_{n=1}^N \log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2),$$

$$\log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2) = -\frac{1}{2\sigma^2} (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const}$$

With

$$\log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2) = -\frac{1}{2\sigma^2} (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const}$$

we get

$$\log p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) = \sum_{n=1}^N \log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const}$$

Maximum Likelihood (2)

With

$$\log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2) = -\frac{1}{2\sigma^2} (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const}$$

we get

$$\begin{aligned}\log p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) &= \sum_{n=1}^N \log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const} \\ &= -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \text{const}\end{aligned}$$

Maximum Likelihood (2)

With

$$\log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2) = -\frac{1}{2\sigma^2} (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const}$$

we get

$$\begin{aligned}\log p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) &= \sum_{n=1}^N \log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const} \\ &= -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \text{const} \\ &= -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2 + \text{const}\end{aligned}$$

Maximum Likelihood (2)

With

$$\log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2) = -\frac{1}{2\sigma^2} (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const}$$

we get

$$\begin{aligned}\log p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) &= \sum_{n=1}^N \log \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const} \\ &= -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \text{const} \\ &= -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2 + \text{const}\end{aligned}$$

- Computing the gradient with respect to $\boldsymbol{\theta}$ and setting it to 0 gives the **maximum likelihood estimator** (least-squares estimator)

$$\boldsymbol{\theta}^{\text{ML}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

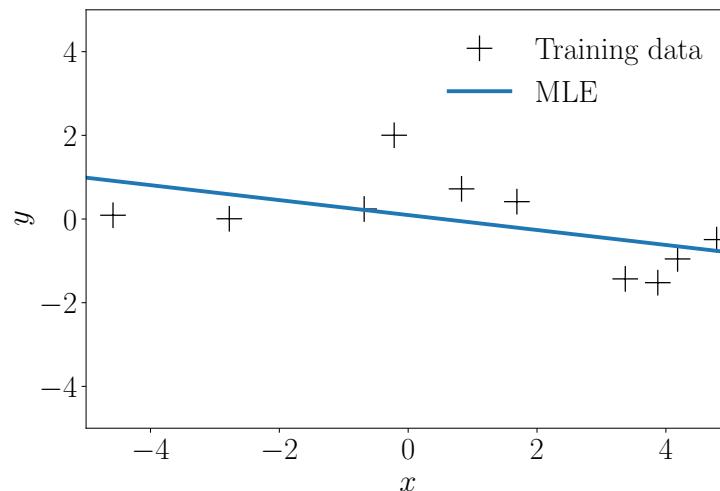
$$y = \mathbf{x}^\top \boldsymbol{\theta} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

Given an arbitrary input \mathbf{x}_* , we can predict the corresponding observation y_* using the maximum likelihood parameter:

$$p(y_* | \mathbf{x}_*, \boldsymbol{\theta}^{\text{ML}}) = \mathcal{N}(y_* | \mathbf{x}_*^\top \boldsymbol{\theta}^{\text{ML}}, \sigma^2)$$

- Measurement noise variance σ^2 assumed known
- In the absence of noise ($\sigma^2 = 0$), the prediction will be deterministic

Example 1: Linear Functions



$$y = \theta_0 + \theta_1 x + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

- At any query point x_* we obtain the mean prediction as

$$\mathbb{E}[y_* | \boldsymbol{\theta}^{\text{ML}}, x_*] = \theta_0^{\text{ML}} + \theta_1^{\text{ML}} x_*$$

$$y = \phi(\mathbf{x})^\top \boldsymbol{\theta} + \epsilon = \sum_{m=0}^M \theta_m x^m + \epsilon$$

- Polynomial regression with features

$$\phi(\mathbf{x}) = [1, x, x^2, \dots, x^M]^\top$$

- Maximum likelihood estimator:

$$y = \phi(x)^\top \theta + \epsilon = \sum_{m=0}^M \theta_m x^m + \epsilon$$

- Polynomial regression with features

$$\phi(x) = [1, x, x^2, \dots, x^M]^\top$$

- Maximum likelihood estimator:

$$\theta^{\text{ML}} = (\Phi^\top \Phi)^{-1} \Phi^\top y$$

Example 2: Polynomial Regression

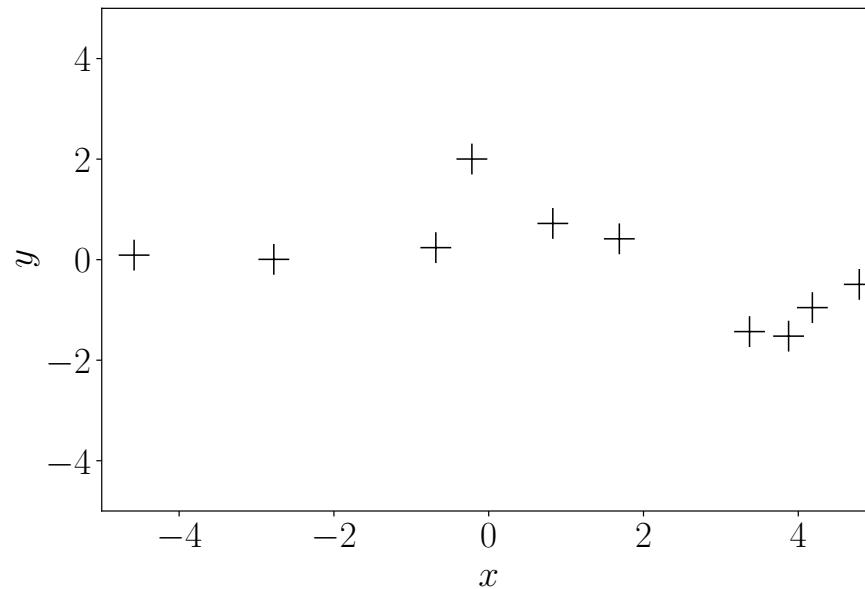


Figure: Training data

Example 2: Polynomial Regression

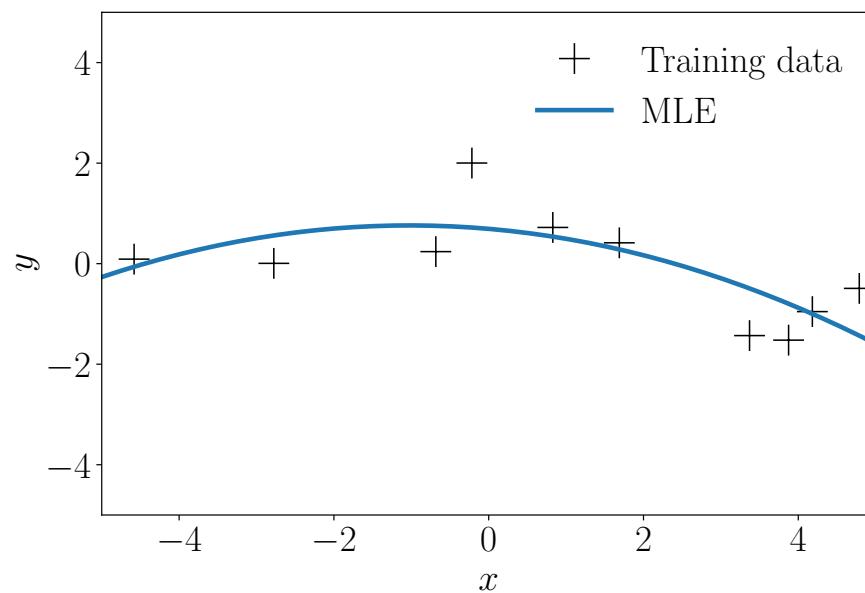


Figure: 2nd-order polynomial

Example 2: Polynomial Regression

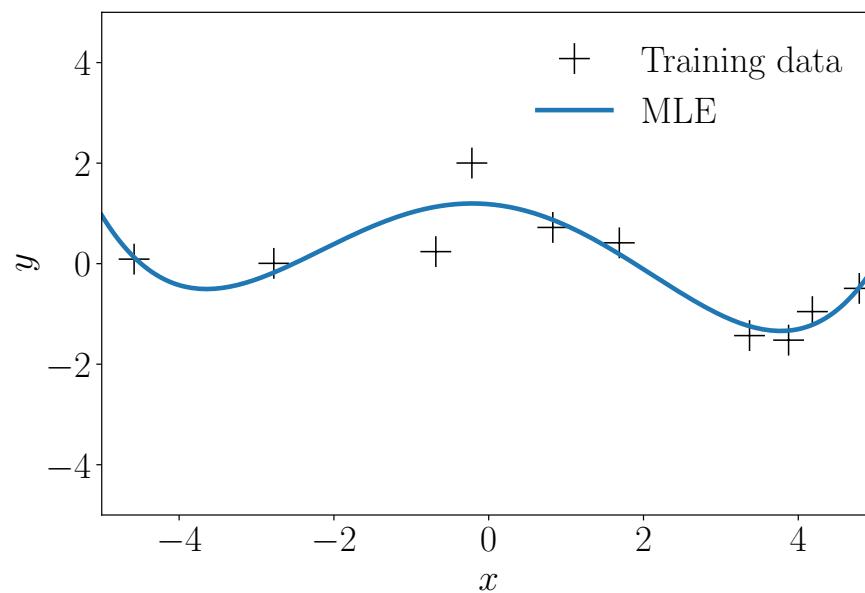


Figure: 4th-order polynomial

Example 2: Polynomial Regression

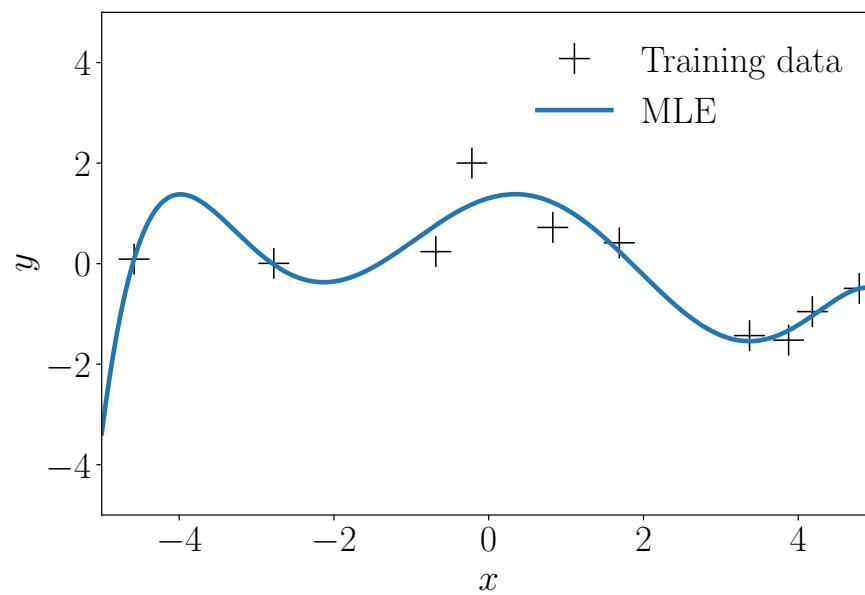


Figure: 6th-order polynomial

Example 2: Polynomial Regression

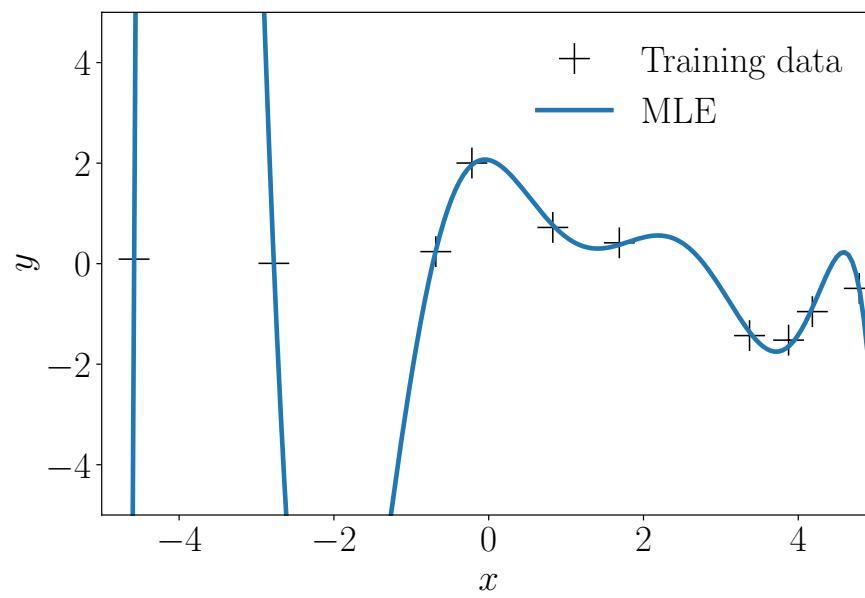


Figure: 8th-order polynomial

Example 2: Polynomial Regression

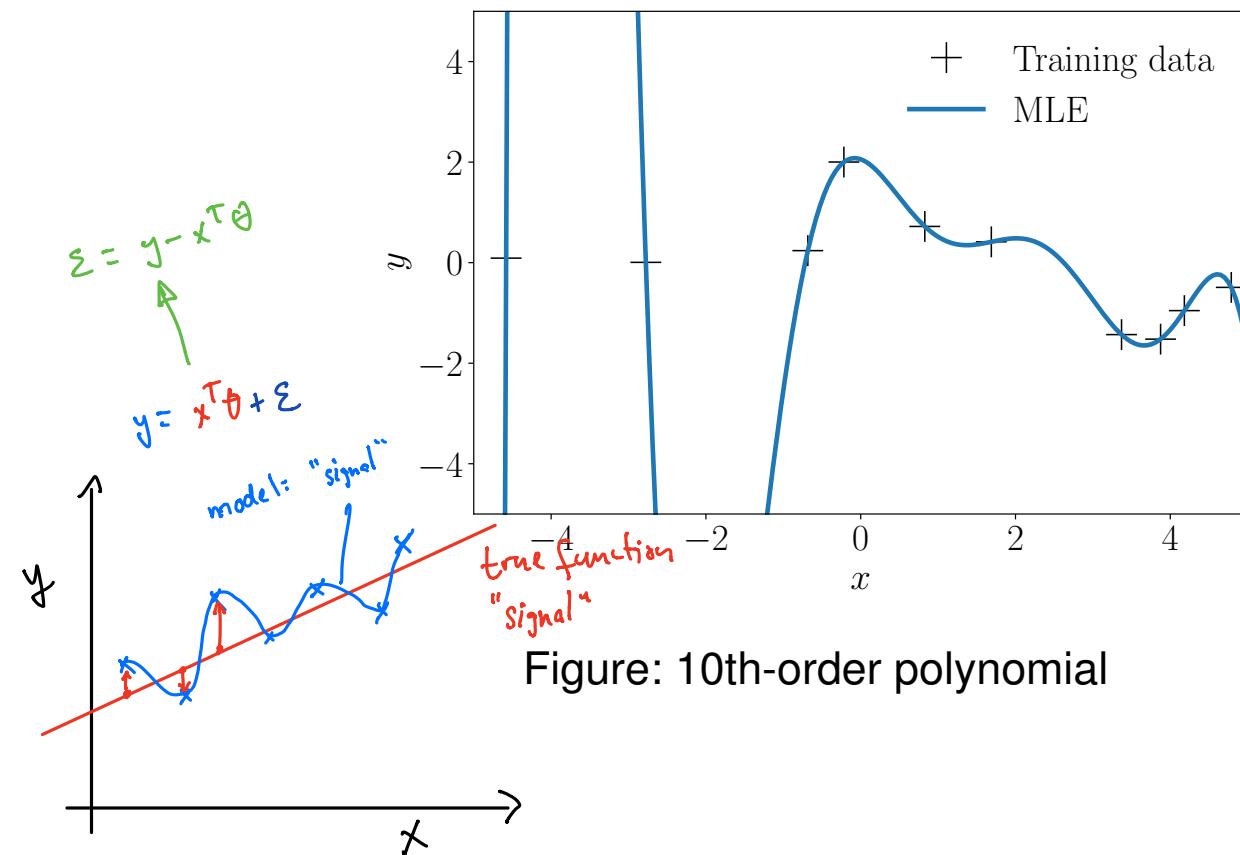
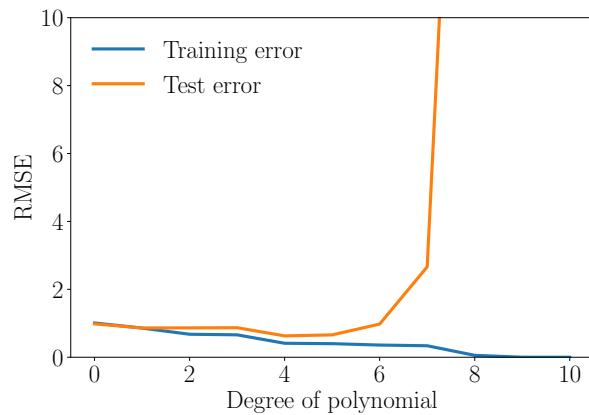


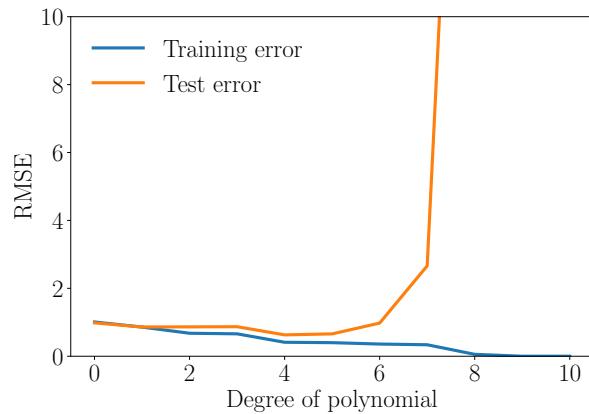
Figure: 10th-order polynomial

Overfitting



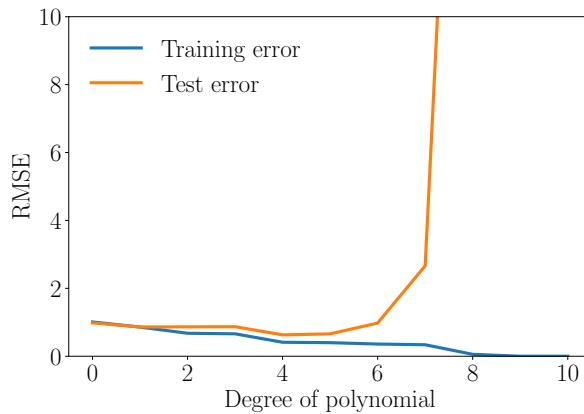
- Training error decreases with higher flexibility of the model

Overfitting



- Training error decreases with higher flexibility of the model
- We are not so much interested in the training error, but in the **generalization error**: How well does the model perform when we predict at previously unseen input locations?

Overfitting



- Training error decreases with higher flexibility of the model
- We are not so much interested in the training error, but in the **generalization error**: How well does the model perform when we predict at previously unseen input locations?
- Maximum likelihood often runs into **overfitting** problems, i.e., we exploit the flexibility of the model to fit to the noise in the data

- **Empirical observation:** Parametric models that overfit tend to have some extreme (large amplitude) parameter values

- **Empirical observation:** Parametric models that overfit tend to have some extreme (large amplitude) parameter values
- Mitigate the effect of overfitting by placing a prior distribution $p(\theta)$ on the parameters
 - ▶▶ Penalize extreme values that are implausible under that prior

- **Empirical observation:** Parametric models that overfit tend to have some extreme (large amplitude) parameter values
- Mitigate the effect of overfitting by placing a prior distribution $p(\theta)$ on the parameters
 - ▶ Penalize extreme values that are implausible under that prior
- Choose θ^* as the parameter that maximizes the (log) parameter posterior

$$\log p(\theta|X, y) = \underbrace{\log p(y|X, \theta)}_{\text{log-likelihood}} + \underbrace{\log p(\theta)}_{\text{log-prior}} + \text{const}$$

- **Empirical observation:** Parametric models that overfit tend to have some extreme (large amplitude) parameter values
- Mitigate the effect of overfitting by placing a prior distribution $p(\theta)$ on the parameters
 - ▶ Penalize extreme values that are implausible under that prior
- Choose θ^* as the parameter that maximizes the (log) parameter posterior

$$\log p(\theta|X, y) = \underbrace{\log p(y|X, \theta)}_{\text{log-likelihood}} + \underbrace{\log p(\theta)}_{\text{log-prior}} + \text{const}$$

- Log-prior induces a direct penalty on the parameters

- **Empirical observation:** Parametric models that overfit tend to have some extreme (large amplitude) parameter values
- Mitigate the effect of overfitting by placing a prior distribution $p(\theta)$ on the parameters
 - ▶ Penalize extreme values that are implausible under that prior
- Choose θ^* as the parameter that maximizes the (log) parameter posterior

$$\log p(\theta|X, y) = \underbrace{\log p(y|X, \theta)}_{\text{log-likelihood}} + \underbrace{\log p(\theta)}_{\text{log-prior}} + \text{const}$$

- Log-prior induces a direct penalty on the parameters
- **Maximum a posteriori estimate** (regularized least squares)

MAP Estimation (2)

$$\nearrow \propto \exp\left(-\frac{1}{2} \theta^\top \alpha^{-2} \Sigma \theta\right) = \exp\left(-\frac{1}{2\alpha^2} \theta^\top \theta\right)$$

- Gaussian parameter prior $p(\theta) = \mathcal{N}(\mathbf{0}, \alpha^2 \mathbf{I})$

- Log-posterior distribution:

$$\log p(\theta | X, y) = -\frac{1}{2\sigma^2} (y - X\theta)^\top (y - X\theta) - \frac{1}{2\alpha^2} \theta^\top \theta + \text{const}$$

$$= -\frac{1}{2\sigma^2} \|y - X\theta\|^2 - \frac{1}{2\alpha^2} \|\theta\|^2 + \text{const}$$

- Gaussian parameter prior $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \alpha^2 \mathbf{I})$
- Log-posterior distribution:

$$\begin{aligned}\log p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}) &= -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) - \frac{1}{2\alpha^2} \boldsymbol{\theta}^\top \boldsymbol{\theta} + \text{const} \\ &= -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2 - \frac{1}{2\alpha^2} \|\boldsymbol{\theta}\|^2 + \text{const}\end{aligned}$$

- Compute gradient with respect to $\boldsymbol{\theta}$, set it to $\mathbf{0}$
► **Maximum a posteriori estimate:**

$$\boldsymbol{\theta}^{\text{MAP}} = (\mathbf{X}^\top \mathbf{X} + \frac{\sigma^2}{\alpha^2} \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

Example: Polynomial Regression

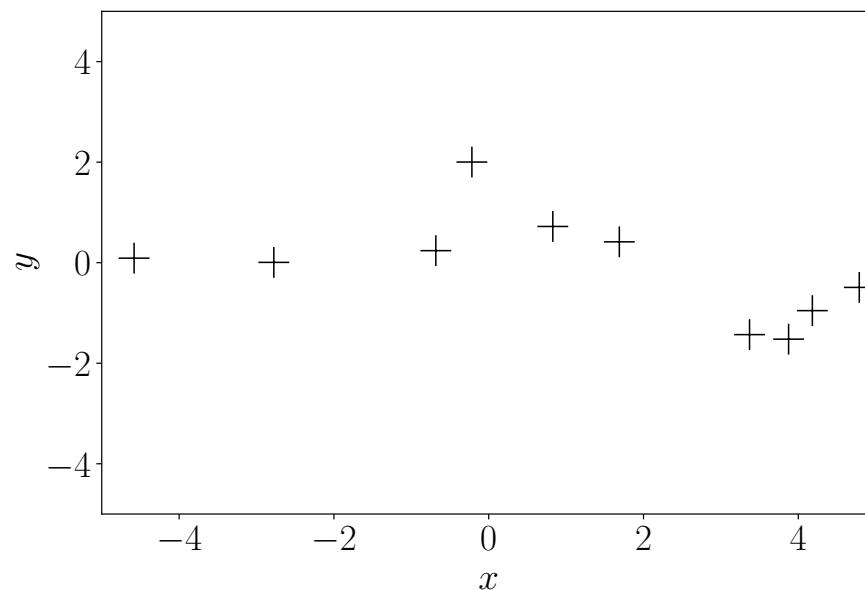


Figure: Training data

Mean prediction:

$$\mathbb{E}[y_* | x_*, \theta^{\text{MAP}}] = \phi^\top(x_*) \theta^{\text{MAP}}$$

Example: Polynomial Regression

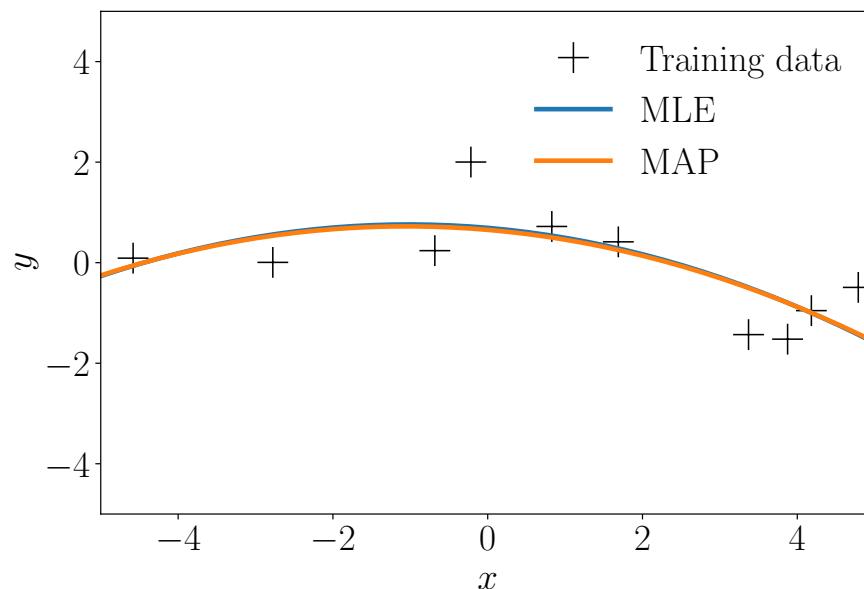


Figure: 2nd-order polynomial

Mean prediction:

$$\mathbb{E}[y_* | \mathbf{x}_*, \boldsymbol{\theta}^{\text{MAP}}] = \boldsymbol{\phi}^\top(\mathbf{x}_*) \boldsymbol{\theta}^{\text{MAP}}$$

Example: Polynomial Regression

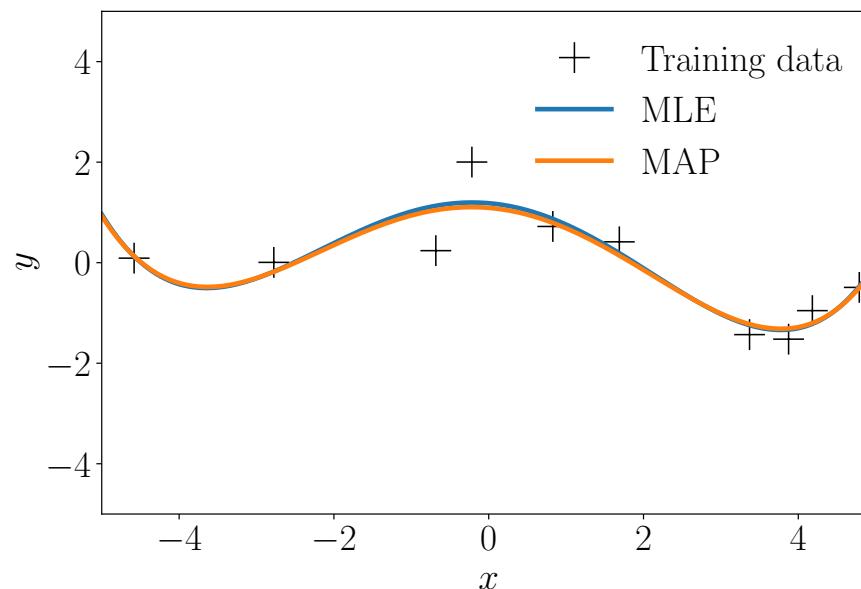


Figure: 4th-order polynomial

Mean prediction:

$$\mathbb{E}[y_* | \mathbf{x}_*, \boldsymbol{\theta}^{\text{MAP}}] = \boldsymbol{\phi}^\top(\mathbf{x}_*) \boldsymbol{\theta}^{\text{MAP}}$$

Example: Polynomial Regression

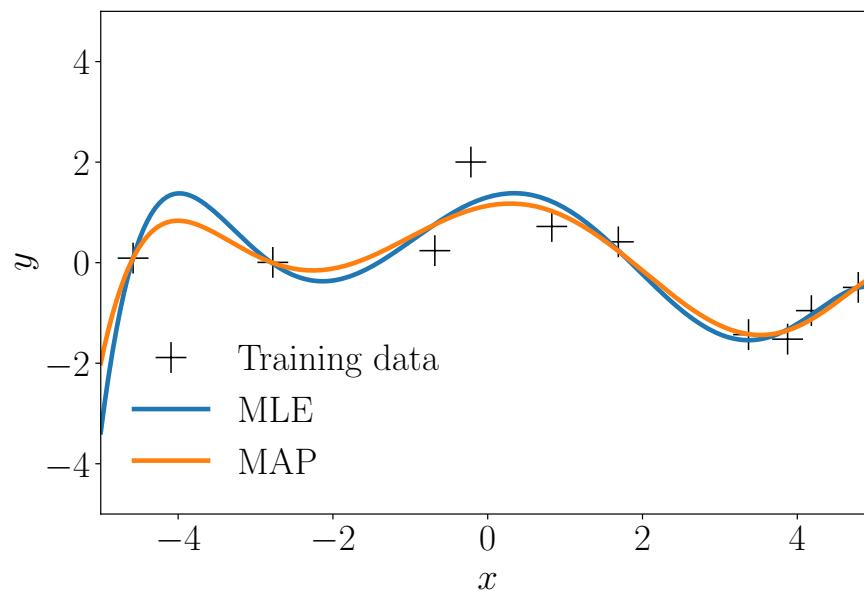


Figure: 6th-order polynomial

Mean prediction:

$$\mathbb{E}[y_* | \mathbf{x}_*, \boldsymbol{\theta}^{\text{MAP}}] = \boldsymbol{\phi}^\top(\mathbf{x}_*) \boldsymbol{\theta}^{\text{MAP}}$$

Example: Polynomial Regression

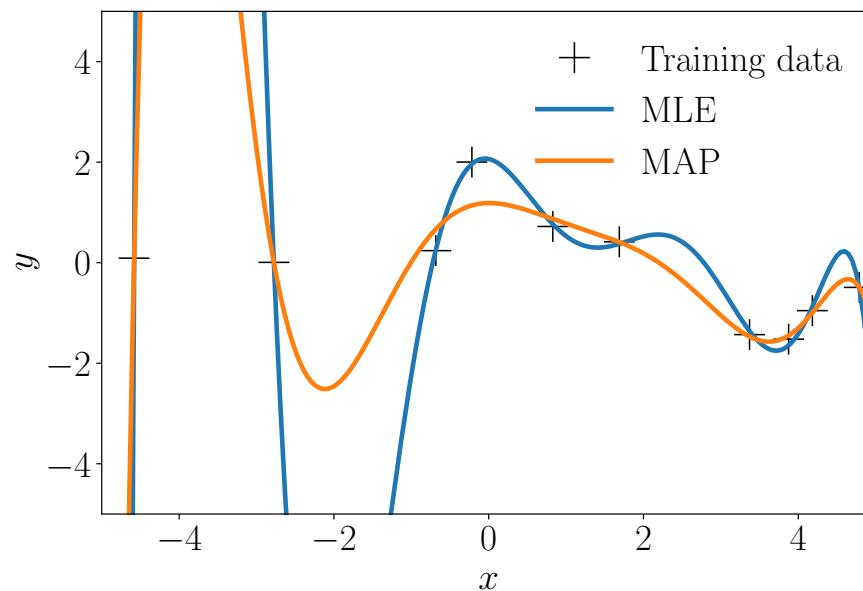


Figure: 8th-order polynomial

Mean prediction:

$$\mathbb{E}[y_* | \mathbf{x}_*, \boldsymbol{\theta}^{\text{MAP}}] = \boldsymbol{\phi}^\top(\mathbf{x}_*) \boldsymbol{\theta}^{\text{MAP}}$$

Example: Polynomial Regression

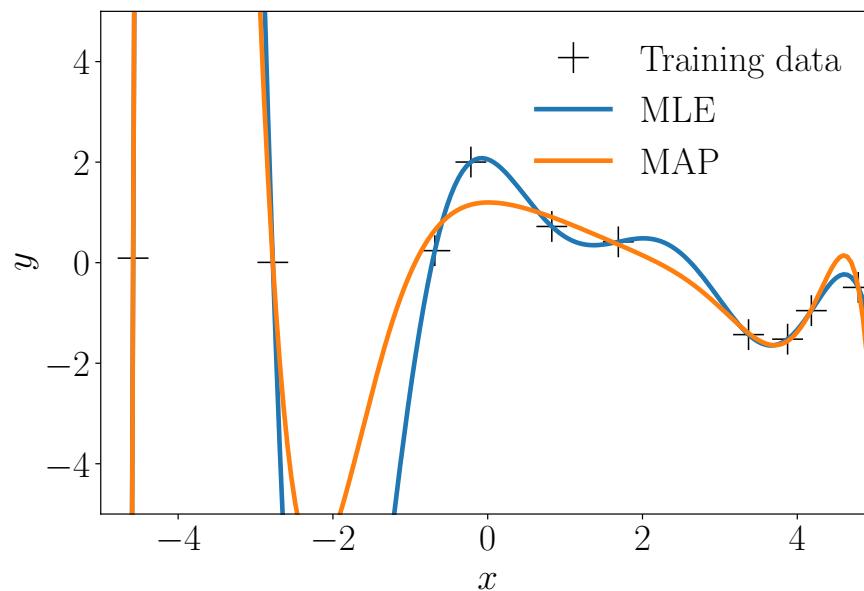
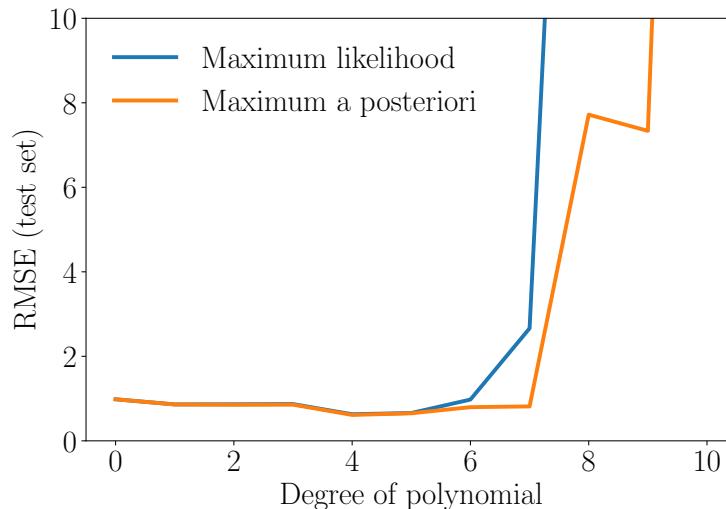


Figure: 10th-order polynomial

Mean prediction:

$$\mathbb{E}[y_* | \mathbf{x}_*, \boldsymbol{\theta}^{\text{MAP}}] = \boldsymbol{\phi}^\top(\mathbf{x}_*) \boldsymbol{\theta}^{\text{MAP}}$$

Generalization Error



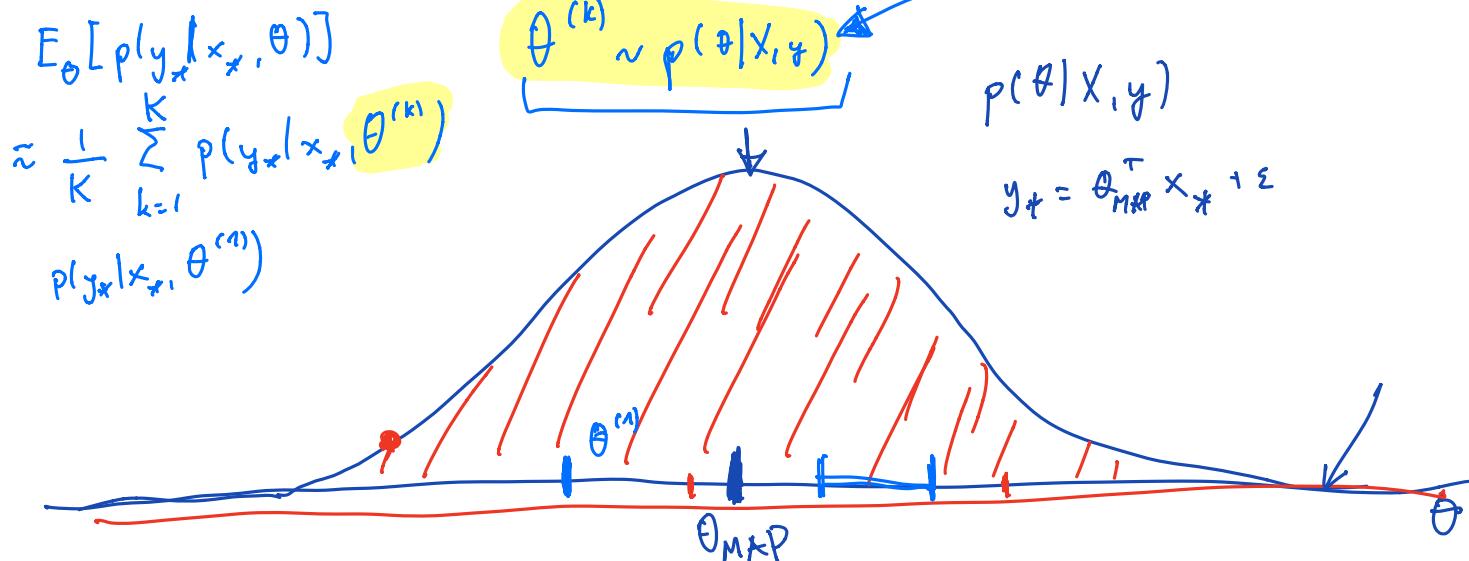
- MAP estimation “delays” the problem of overfitting
- It does not provide a general solution
- ▶ Need a more principled solution

Bayesian Linear Regression

$$y = \phi^\top(\mathbf{x})\theta + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

- Avoid overfitting by not fitting any parameters:

► Integrate parameters out instead of optimizing them



Bayesian Linear Regression

$$y = \phi^\top(\mathbf{x})\theta + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

- Avoid overfitting by not fitting any parameters:

► Integrate parameters out instead of optimizing them

- Use a full parameter distribution $p(\theta)$ (and not a single point estimate θ^*) when making predictions:

$$p(y_*|x_*) = \int p(y_*|x_*, \theta) p(\theta) d\theta$$

plausible values
of θ

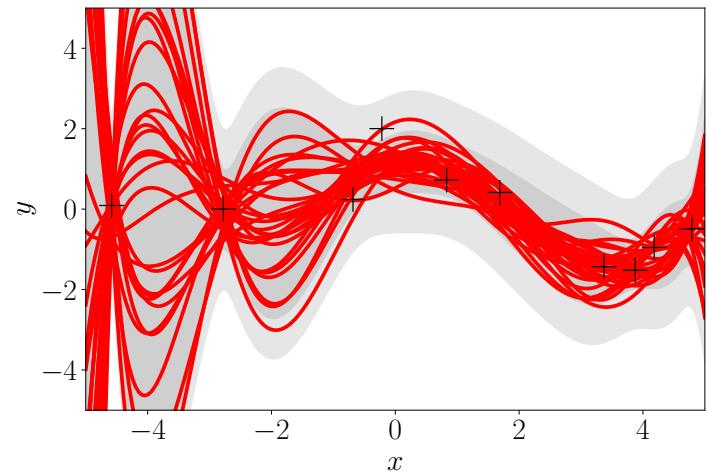
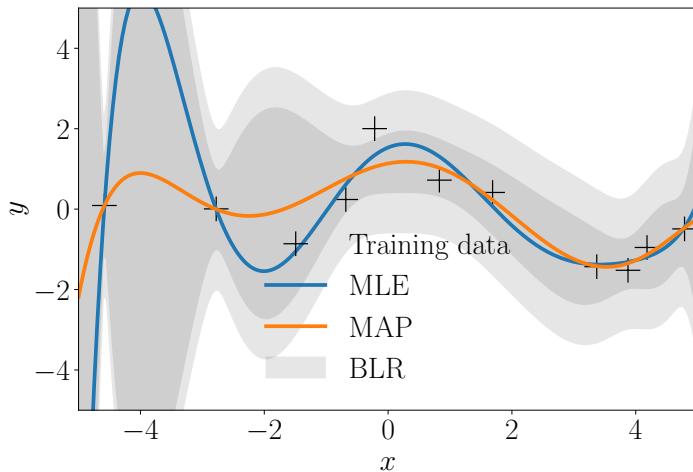
previously:
 $p(y_*|x_*, \theta_{MAP/ML})$

$\approx \underbrace{E_{\theta} [p(y_*|x_*, \theta)]}$

pred. distribution
given a single param.
value θ

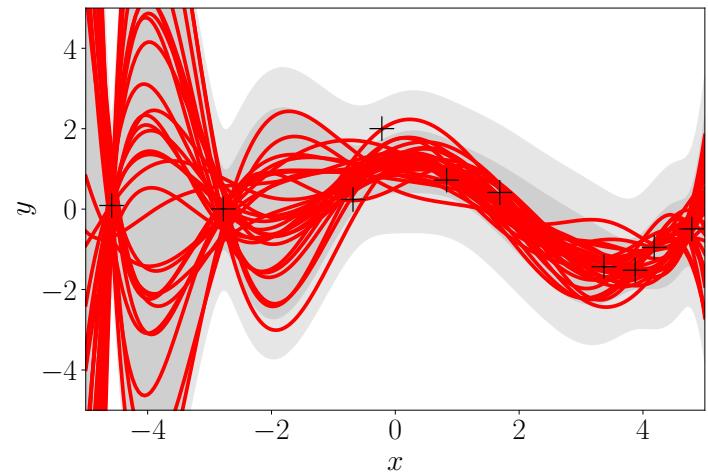
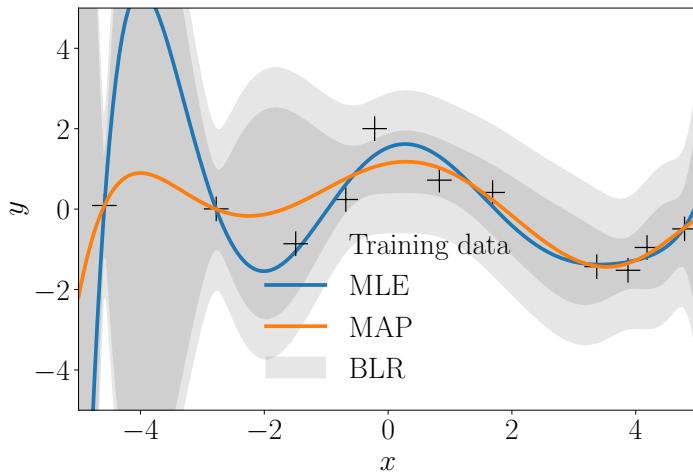
- Prediction no longer depends on θ
- Predictive distribution reflects the uncertainty about the “correct” parameter setting

Example



- Light-gray: uncertainty due to noise (same as in MLE/MAP)
- Dark-gray: uncertainty due to parameter uncertainty

Example



- Light-gray: uncertainty due to noise (same as in MLE/MAP)
- Dark-gray: uncertainty due to parameter uncertainty
- Right: Plausible functions under the parameter distribution (every single parameter setting describes one function)

Model for Bayesian Linear Regression

$$\exp\left(-\frac{1}{2\sigma^2}(y - \phi^\top \theta)^2\right) = \dots \exp\left(-\frac{1}{2}(\theta - m)^\top S^{-1}(\theta - m)\right)$$

$-\frac{1}{2\sigma^2} \theta^\top \theta + \text{const.}$

Prior $p(\theta) = \mathcal{N}(\boldsymbol{m}_0, \boldsymbol{S}_0),$

Likelihood $p(y|\boldsymbol{x}, \theta) = \mathcal{N}(y | \boldsymbol{\phi}^\top(\boldsymbol{x})\theta, \sigma^2)$

- Parameter θ becomes a latent (random) variable
- Prior distribution induces a **distribution over plausible functions**
- Choose a **conjugate Gaussian prior**

- Closed-form computations
- Gaussian posterior

posterior \propto prior \times likelihood

Given: likelihood

Choose: prior such that the posterior has the same functional form as the prior

e.g.: Gauss. likelihood $\mathcal{N}(y | \boldsymbol{x}^\top \theta, \sigma^2)$

choose: Gauss. prior $\mathcal{N}(\theta | \boldsymbol{0}, \boldsymbol{\alpha}^{-1})$

\Rightarrow posterior: Gaussian $\mathcal{N}(\theta | \boldsymbol{m}, \boldsymbol{S})$

Parameter Posterior and Predictions

■ Prior $p(\theta) = \mathcal{N}(m_0, S_0)$ is Gaussian \rightarrow posterior is Gaussian:

\rightarrow Derive this

$$\text{var}_{\theta}[A\theta] = A \text{var}[\theta] A^T$$

$$p(\theta|X, y) = \mathcal{N}(m_N, S_N)$$

$$S_N = (S_0^{-1} + \sigma^{-2} \Phi^\top \Phi)^{-1}$$

$$m_N = S_N(S_0^{-1}m_0 + \sigma^{-2} \Phi^\top y)$$

$$p(\theta|x, y) = \mathcal{N}(m_N, S_N)$$

$$y_* = \phi^\top(x_*) \theta + \varepsilon$$

$$E[y_*|x_*] = E[\phi^\top(x_*) \theta + \varepsilon] = \underbrace{\phi^\top(x_*)}_{\text{const}} \underbrace{E[\theta]}_{\hat{\theta} = m_N} + \cancel{E[\varepsilon]} = \phi^\top(x_*) m_N$$

$$\text{var}[y_*|x_*] = \text{var}[\phi^\top(x_*) \theta + \varepsilon] = \text{var}[\phi^\top(x_*) \theta] + \text{var}[\varepsilon]$$

$$= \phi^\top(x_*) \underbrace{\text{var}[\theta]}_{S_N} \phi(x_*) + \sigma^2 = \phi^\top(x_*) S_N \phi(x_*) + \sigma^2$$

Parameter Posterior and Predictions

- Prior $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{m}_0, \boldsymbol{S}_0)$ is Gaussian ➤ posterior is Gaussian:

$$p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}) = \mathcal{N}(\boldsymbol{m}_N, \boldsymbol{S}_N)$$

$$\boldsymbol{S}_N = (\boldsymbol{S}_0^{-1} + \sigma^{-2} \boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1}$$

$$\boldsymbol{m}_N = \boldsymbol{S}_N (\boldsymbol{S}_0^{-1} \boldsymbol{m}_0 + \sigma^{-2} \boldsymbol{\Phi}^\top \mathbf{y})$$

- Mean \boldsymbol{m}_N identical to MAP estimate

Parameter Posterior and Predictions

- Prior $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{m}_0, \boldsymbol{S}_0)$ is Gaussian \rightarrow posterior is Gaussian:

$$p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}) = \mathcal{N}(\boldsymbol{m}_N, \boldsymbol{S}_N)$$

$$\boldsymbol{S}_N = (\boldsymbol{S}_0^{-1} + \sigma^{-2} \boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1}$$

$$\boldsymbol{m}_N = \boldsymbol{S}_N (\boldsymbol{S}_0^{-1} \boldsymbol{m}_0 + \sigma^{-2} \boldsymbol{\Phi}^\top \mathbf{y})$$

- Mean \boldsymbol{m}_N identical to MAP estimate
- Assume a Gaussian distribution $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{m}_N, \boldsymbol{S}_N)$. Then

$$p(y_* | \mathbf{x}_*) = \mathcal{N}(y | \boldsymbol{\phi}^\top(\mathbf{x}_*) \boldsymbol{m}_N, \boldsymbol{\phi}^\top(\mathbf{x}_*) \boldsymbol{S}_N \boldsymbol{\phi}(\mathbf{x}_*) + \sigma^2)$$

Parameter Posterior and Predictions

- Prior $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{m}_0, \boldsymbol{S}_0)$ is Gaussian ➤ posterior is Gaussian:

$$p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}) = \mathcal{N}(\boldsymbol{m}_N, \boldsymbol{S}_N)$$

$$\boldsymbol{S}_N = (\boldsymbol{S}_0^{-1} + \sigma^{-2} \boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1}$$

$$\boldsymbol{m}_N = \boldsymbol{S}_N (\boldsymbol{S}_0^{-1} \boldsymbol{m}_0 + \sigma^{-2} \boldsymbol{\Phi}^\top \mathbf{y})$$

- Mean \boldsymbol{m}_N identical to MAP estimate
- Assume a Gaussian distribution $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{m}_N, \boldsymbol{S}_N)$. Then

$$p(y_* | \mathbf{x}_*) = \mathcal{N}(y | \boldsymbol{\phi}^\top(\mathbf{x}_*) \boldsymbol{m}_N, \boldsymbol{\phi}^\top(\mathbf{x}_*) \boldsymbol{S}_N \boldsymbol{\phi}(\mathbf{x}_*) + \sigma^2)$$

- $\boldsymbol{\phi}^\top(\mathbf{x}_*) \boldsymbol{S}_N \boldsymbol{\phi}(\mathbf{x}_*)$: Accounts for parameter uncertainty in predictive variance

More details ➤ <https://mml-book.com>, Chapter 9

- Marginal likelihood can be computed analytically.
- With $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

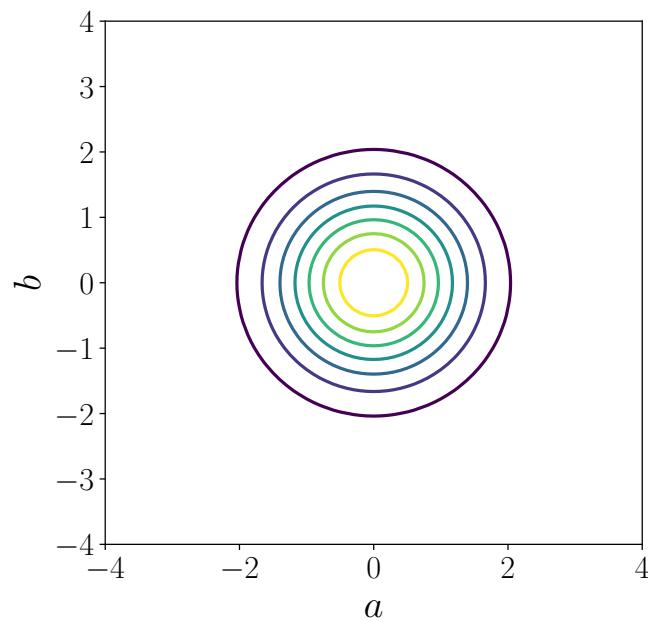
$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta} = \mathcal{N}(\mathbf{y} | \boldsymbol{\Phi}\boldsymbol{\mu}, \boldsymbol{\Phi}\boldsymbol{\Sigma}\boldsymbol{\Phi}^\top + \sigma^2 \mathbf{I})$$

- Derivation via completing the squares (see Section 9.3.5 of MML book)

Distribution over Functions

Consider a linear regression setting

$$y = f(x) + \epsilon = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$
$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$



Sampling from the Prior over Functions

Consider a linear regression setting

$$y = f(x) + \epsilon = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$

$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$f_i(x) = a_i + b_i x, \quad [a_i, b_i] \sim p(a, b)$$

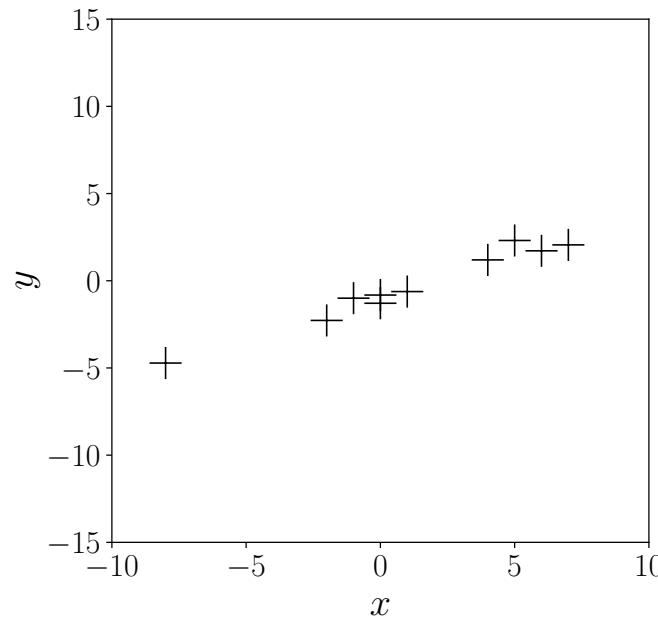
Sampling from the Posterior over Functions

Consider a linear regression setting

$$y = f(x) + \epsilon = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$

$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$\mathbf{X} = [x_1, \dots, x_N]$, $\mathbf{y} = [y_1, \dots, y_N]$ Training inputs/targets



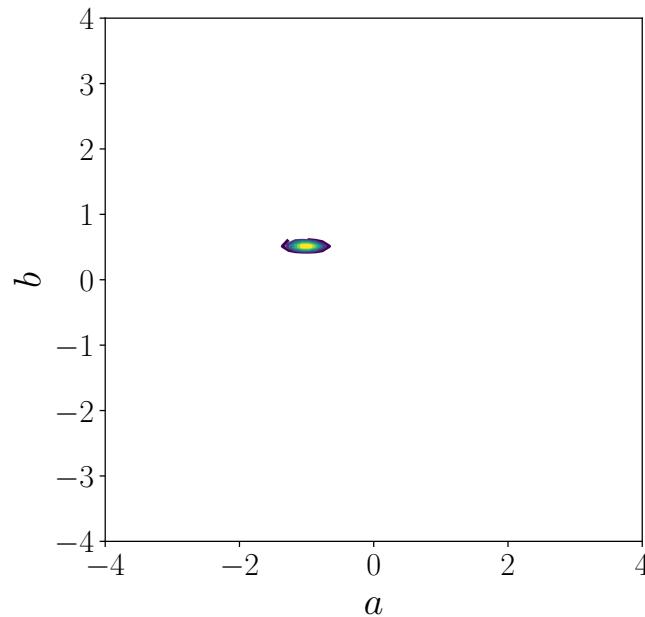
Sampling from the Posterior over Functions

Consider a linear regression setting

$$y = f(x) + \epsilon = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$

$$p(a, b) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$p(a, b | \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{m}_N, \mathbf{S}_N) \quad \text{Posterior}$$



Sampling from the Posterior over Functions



Consider a linear regression setting

$$y = f(x) + \epsilon = a + bx + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$

$$[a_i, b_i] \sim p(a, b | \mathbf{X}, \mathbf{y})$$

$$f_i = a_i + b_i x$$

Fitting Nonlinear Functions

- Fit nonlinear functions using (Bayesian) linear regression:
Linear combination of nonlinear features

Fitting Nonlinear Functions

- Fit nonlinear functions using (Bayesian) linear regression:
Linear combination of nonlinear features
- Example: Radial-basis-function (RBF) network

$$f(\boldsymbol{x}) = \sum_{i=1}^n \theta_i \phi_i(\boldsymbol{x}), \quad \theta_i \sim \mathcal{N}(0, \sigma_p^2)$$

Fitting Nonlinear Functions

- Fit nonlinear functions using (Bayesian) linear regression:
Linear combination of nonlinear features
- Example: Radial-basis-function (RBF) network

$$f(\boldsymbol{x}) = \sum_{i=1}^n \theta_i \phi_i(\boldsymbol{x}), \quad \theta_i \sim \mathcal{N}(0, \sigma_p^2)$$

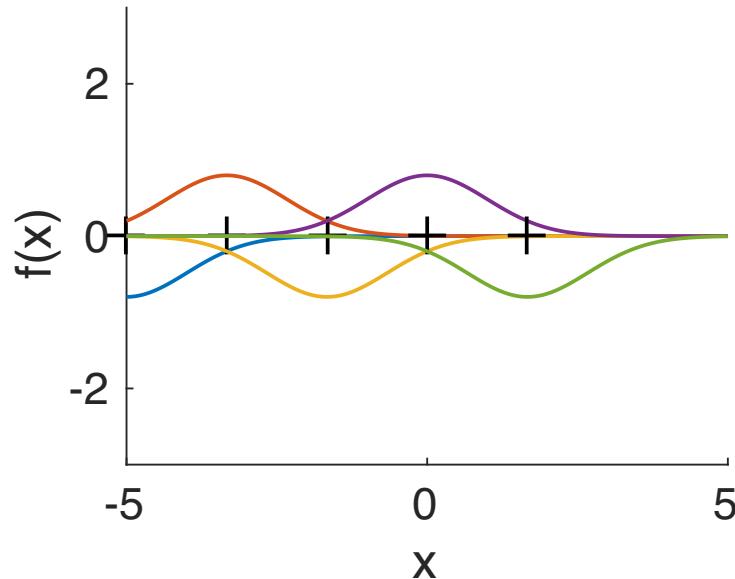
where

$$\phi_i(\boldsymbol{x}) = \exp\left(-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu}_i)^\top(\boldsymbol{x} - \boldsymbol{\mu}_i)\right)$$

for given “centers” $\boldsymbol{\mu}_i$

Illustration: Fitting a Radial Basis Function Network

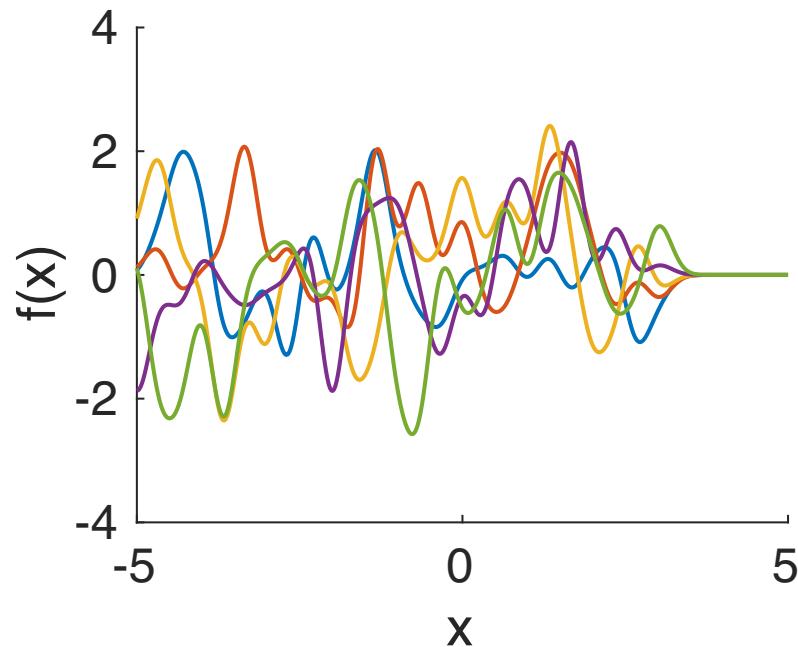
$$\phi_i(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^\top(\mathbf{x} - \boldsymbol{\mu}_i)\right)$$



- Place Gaussian-shaped basis functions ϕ_i at 25 input locations μ_i , linearly spaced in the interval $[-5, 3]$

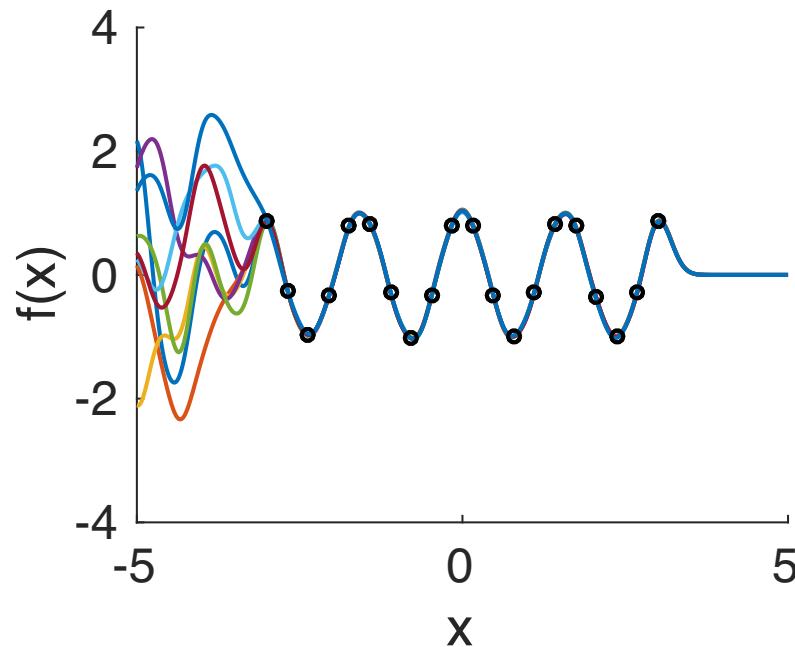
Samples from the RBF Prior

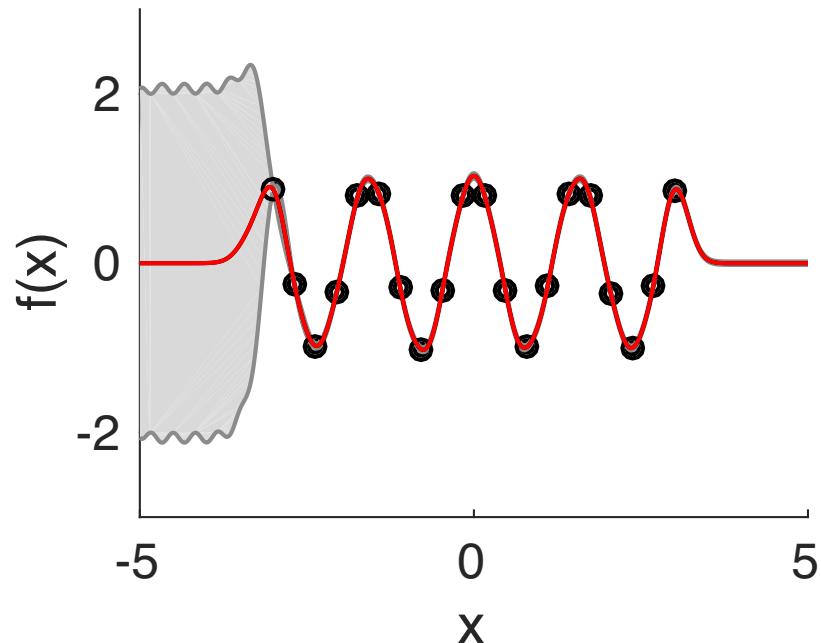
$$f(\boldsymbol{x}) = \sum_{i=1}^n \theta_i \phi_i(\boldsymbol{x}), \quad p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

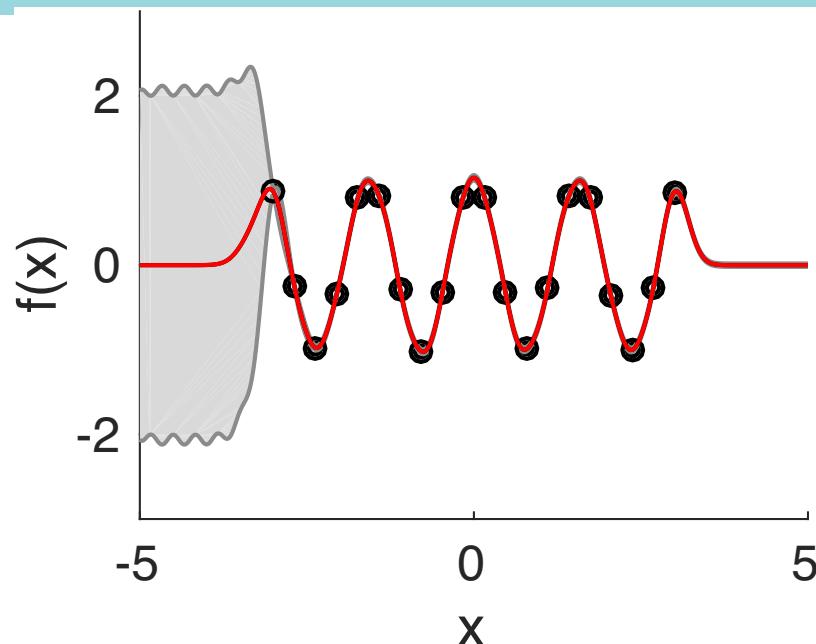


Samples from the RBF Posterior

$$f(\boldsymbol{x}) = \sum_{i=1}^n \theta_i \phi_i(\boldsymbol{x}), \quad p(\boldsymbol{\theta} | \boldsymbol{X}, \boldsymbol{y}) = \mathcal{N}(\boldsymbol{m}_N, \boldsymbol{S}_N)$$







- Feature engineering (what basis functions to use?)
- Finite number of features:
 - Above: Without basis functions on the right, we cannot express any variability of the function
 - Ideally: Add more (infinitely many) basis functions

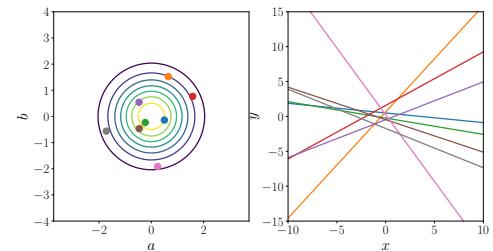
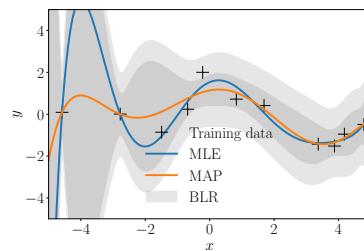
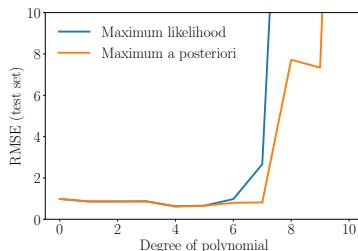
- Instead of sampling parameters, which induce a distribution over functions, **sample functions directly**
 - ▶ Place a prior on functions
 - ▶ Make assumptions on the distribution of functions

- Instead of sampling parameters, which induce a distribution over functions, **sample functions directly**
 - ▶ Place a prior on functions
 - ▶ Make assumptions on the distribution of functions
- Intuition: function = infinitely long vector of function values
 - ▶ Make assumptions on the distribution of function values

- Instead of sampling parameters, which induce a distribution over functions, **sample functions directly**
 - ▶ Place a prior on functions
 - ▶ Make assumptions on the distribution of functions
- Intuition: function = infinitely long vector of function values
 - ▶ Make assumptions on the distribution of function values

- Instead of sampling parameters, which induce a distribution over functions, **sample functions directly**
 - ▶ Place a prior on functions
 - ▶ Make assumptions on the distribution of functions
 - Intuition: function = infinitely long vector of function values
 - ▶ Make assumptions on the distribution of function values
- ▶ **Gaussian process**

Summary



- Regression = curve fitting
- Linear regression = linear in the parameters
- Parameter estimation via maximum likelihood and MAP estimation can lead to **overfitting**
- **Bayesian linear regression** addresses this issue, but may not be analytically tractable
- Predictive uncertainty in Bayesian linear regression explicitly accounts for parameter uncertainty
- Distribution over parameters ➤ Distribution over functions

Appendix

■ Joint Gaussian distribution

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{bmatrix} \right)$$

■ Joint Gaussian distribution

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{bmatrix} \right)$$

■ Marginal:

$$\begin{aligned} p(\mathbf{x}) &= \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y} \\ &= \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_{xx}) \end{aligned}$$

Joint Gaussian Distribution

■ Joint Gaussian distribution

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{bmatrix} \right)$$

■ Marginal:

$$\begin{aligned} p(\mathbf{x}) &= \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y} \\ &= \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_{xx}) \end{aligned}$$

■ Conditional:

$$\begin{aligned} p(\mathbf{x}|\mathbf{y}) &= \mathcal{N}(\boldsymbol{\mu}_{x|y}, \boldsymbol{\Sigma}_{x|y}) \\ \boldsymbol{\mu}_{x|y} &= \boldsymbol{\mu}_x + \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} (\mathbf{y} - \boldsymbol{\mu}_y) \\ \boldsymbol{\Sigma}_{x|y} &= \boldsymbol{\Sigma}_{xx} - \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yx} \end{aligned}$$

If $x \sim \mathcal{N}(x | \mu, \Sigma)$ and $z = Ax + b$ then

$$p(z) = \mathcal{N}(z | A\mu + b, A\Sigma A^\top)$$

Product of Two Gaussians

$x \in \mathbb{R}^D$. Then:

$$\mathcal{N}(x | \mathbf{a}, \mathbf{A})\mathcal{N}(x | \mathbf{b}, \mathbf{B}) = Z\mathcal{N}(x | \mathbf{c}, \mathbf{C})$$

$$\mathbf{C} = (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1}$$

$$\mathbf{c} = \mathbf{C}(\mathbf{A}^{-1}\mathbf{a} + \mathbf{B}^{-1}\mathbf{b})$$

$$Z = (2\pi)^{-\frac{D}{2}} |\mathbf{A} + \mathbf{B}| \exp\left(-\frac{1}{2}(\mathbf{a} - \mathbf{b})^\top (\mathbf{A} + \mathbf{B})^{-1}(\mathbf{a} - \mathbf{b})\right)$$

Product of Two Gaussians

$x \in \mathbb{R}^D$. Then:

$$\mathcal{N}(x | \mathbf{a}, \mathbf{A})\mathcal{N}(x | \mathbf{b}, \mathbf{B}) = Z\mathcal{N}(x | \mathbf{c}, \mathbf{C})$$

$$\mathbf{C} = (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1}$$

$$\mathbf{c} = \mathbf{C}(\mathbf{A}^{-1}\mathbf{a} + \mathbf{B}^{-1}\mathbf{b})$$

$$Z = (2\pi)^{-\frac{D}{2}} |\mathbf{A} + \mathbf{B}| \exp\left(-\frac{1}{2}(\mathbf{a} - \mathbf{b})^\top (\mathbf{A} + \mathbf{B})^{-1}(\mathbf{a} - \mathbf{b})\right)$$

- Product of two Gaussians is an unnormalized Gaussian

Product of Two Gaussians

$x \in \mathbb{R}^D$. Then:

$$\mathcal{N}(x | \mathbf{a}, \mathbf{A})\mathcal{N}(x | \mathbf{b}, \mathbf{B}) = Z\mathcal{N}(x | \mathbf{c}, \mathbf{C})$$

$$\mathbf{C} = (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1}$$

$$\mathbf{c} = \mathbf{C}(\mathbf{A}^{-1}\mathbf{a} + \mathbf{B}^{-1}\mathbf{b})$$

$$Z = (2\pi)^{-\frac{D}{2}} |\mathbf{A} + \mathbf{B}| \exp\left(-\frac{1}{2}(\mathbf{a} - \mathbf{b})^\top (\mathbf{A} + \mathbf{B})^{-1}(\mathbf{a} - \mathbf{b})\right)$$

- Product of two Gaussians is an unnormalized Gaussian
- The “un-normalizer” Z has a Gaussian functional form:

$$Z = \mathcal{N}(\mathbf{a} | \mathbf{b}, \mathbf{A} + \mathbf{B}) = \mathcal{N}(\mathbf{b} | \mathbf{a}, \mathbf{A} + \mathbf{B})$$

Note: This is not a distribution (no random variables)

Example: Marginalization of a Product

$$p_1(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mathbf{a}, \mathbf{A})$$
$$p_2(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mathbf{b}, \mathbf{B})$$

Then

$$\int p_1(\mathbf{x})p_2(\mathbf{x})d\mathbf{x} = \in \mathbb{R}$$

Note: In this context, \mathcal{N} is used to describe the functional relationship between \mathbf{a} , \mathbf{b} . Do not treat \mathbf{a} or \mathbf{b} as random variables—they are both deterministic quantities.

Example: Marginalization of a Product

$$p_1(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mathbf{a}, \mathbf{A})$$
$$p_2(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mathbf{b}, \mathbf{B})$$

Then

$$\int p_1(\mathbf{x})p_2(\mathbf{x})d\mathbf{x} = Z = \mathcal{N}(\mathbf{a} | \mathbf{b}, \mathbf{A} + \mathbf{B}) \in \mathbb{R}$$

Note: In this context, \mathcal{N} is used to describe the functional relationship between \mathbf{a} , \mathbf{b} . Do not treat \mathbf{a} or \mathbf{b} as random variables—they are both deterministic quantities.